

[illegible]

The present invention relates generally to digital video signal processing. More particularly, this invention relates to a technique for scalable buffering in a digital video decoder.

5

15 Experts Group (“MPEG”) and Joint Pictures Experts Group (“JPEG”) standards.

20 is applicable to both types of decoders. However, by way of example, the present

by Sun Microsystems, Inc., Palo Alto, California) have been introduced for use in

general purpose computers. These advanced processor instruction sets optimize the computational aspects of standards-based video decoding. However, there is still a need for a standards-based video decoder design that optimizes the heavily pipelined data-path aspects of the video decoder's underlying memory system. CPU

- 5 performance is directly related to the time the CPU spends in executing a program and the time that the CPU waits for the memory system. By reducing memory system access times, CPU performance can be enhanced.

An optimized decoder design involves several challenges. The design must consider effects that its implementation can have on performance aspects of the
10 decoder's underlying memory system. For example, to facilitate rapid access to data and instructions, a general purpose CPU typically uses a Data Cache ("D-cache") and a separate Instruction Cache ("I-cache"). Cache use is optimized when required data and instructions are located in a respective cache. CPU stall-cycles come primarily from cache misses (a cache miss occurs when the necessary data or instructions are not
15 in a cache). To optimize cache use and increase overall processor performance, a video decoder should design its data-path both in a way that cache misses are minimized and caches are not underutilized.

Most existing cache use optimization schemes are optimized for a particular platform. Thus, a single cache use optimization scheme is not readily ported to
20 different computer architectures.

Designers often need to make difficult cache use design tradeoffs to design a video decoder that is portable across several architectures. These tradeoffs involve balancing cache use factors that cannot all be maximized at the same time. The results of implementing these design tradeoffs are unpredictable and often lead to costly
25 software rewrites to accommodate some new knowledge about the costs and benefits of the design tradeoffs.

In view of the foregoing, it would be highly desirable to provide an improved video decoder with scalable buffers that can be dynamically re-sized to optimally process a video input stream.

[illegible]

5

10

20

30

30

FIGURE 4 is a block diagram of a programmed digital image video decoder with scalable buffering according to an embodiment of the present invention.

The analysis module 145 accumulates a set of cache performance results (cache miss rates) in relation to programmed data buffer sizes. The analysis module 145 contains a report generator module 150 for reporting at least a subset of the cache performance results. The analysis module 145 also includes a buffer size adjuster 152. The buffer size adjuster 152 automatically adjusts the size of data store buffers 160 based upon the cache performance results. The buffer size adjuster 152 may be used to set a static buffer size or to accomplish dynamic buffer size adjustments, as discussed below.

Figure 1 also illustrates data store buffers 160. The data store buffers 160 are video data buffers implemented as either software array structures or regions of dedicated Static Random Access Memory (SRAM) that store video data between processing stages. The term buffer is sometimes used in connection with caches and other physical memories. As used herein, the term buffer refers to a data structure or physical memory to store video data between processing stages of a video decoder. Thus, in a software video decoder, the buffer may be a software array structure or a region of dedicated SRAM. In a hardware video decoder, the buffer will be a physical memory region between processing stages.

The data store buffers 160 are scalable buffers utilized by the video decoding module 135. As discussed below, the size of the data store buffers 160 is adjusted and cache performance is analyzed for different data store buffer sizes. Based upon this analysis, an optimal buffer size is subsequently selected and utilized, as discussed below.

Figure 2 illustrates a video decoding module 135, which may be implemented in software or hardware. An encoded image data bitstream is received at an input node 210 and is presented to a Variable Length Decoder/Inverse Quantization Unit (“VLD/IQ”) 230 for decoding of motion vectors and Discrete Cosine Transform (“DCT”) coefficients, resulting in fixed-length data. The fixed-length data is presented to an Inverse Discrete Cosine Transform Unit (“IDCT”) 240, which determines Displaced Frame Difference (“DFD”) information (a decoded prediction error signal). The fixed-length data is also applied to a Motion Compensator (“MC”) 250. The Motion Compensator 250 uses a reference frame signal on line 220 and a motion vector decoded by the VLD/IQ 230 to generate a motion compensated prediction signal. The motion compensated prediction signal is combined with the displaced frame difference information from the IDCT 240 at mixer 260 to produce the final product, a decoded video signal on line 270.

The operation of the VLD/IQ 230, IDCT 240, MC 250, and mixer 260 are well known in the art. The invention is not directed to the independent operation of these components. Rather, the invention is directed toward the utilization of scalable buffers within a video decoder to optimize the performance of the physical memories (e.g., caches) of the video decoder. This facilitates improved decoder performance. For example, in the case of a decoder operating on a general purpose computer, the performance of the data cache and instruction cache of the CPU are enhanced.

Figure 2 illustrates buffers 272, 274, 276, and 278 positioned between processing stages of the decoder 135. Each data store buffer isolates a functional

unit's source and/or destination data bitstream from another functional unit in the decoding pipeline. For example, buffer 272 isolates the output of VLD/IQ 230 and the input to IDCT 240. As long as a source buffer never overflows, and a destination buffer never underflows, each functional unit (e.g., 230, 240, and 250) can be implemented efficiently.

Figure 3 illustrates a compressed digital video image frame 300 divided into a set of macroblocks 310A-310N. Each macroblock 310 in the digital image frame 300 contains a 16x16 matrix of encoded image data pixels. A macroblock is a standard processing segment in standards-based video decoders. Most prior art video decoders process a single macroblock at a time. In the case of a decoder operating on a general purpose computer, decoding only one macroblock of image data at a time causes the instruction cache 113 of the CPU 110 to be over-utilized, while the data cache 113 of the CPU 110 is under-utilized. The data cache 113 is under-utilized because it is storing only one macroblock of a data. The instruction cache is over-utilized because it must invoke all instructions for processing of the macroblock through all of the functional units (e.g., 230, 240, 250, and 260) of the decoder. In this case, overall video decoder performance could have been improved by making better use of the data cache.

Figure 4a illustrates the relationship between instruction cache miss rates and data store buffer sizes. The instruction cache miss rate is highest when a corresponding data store buffer size is small. The instruction cache miss rate decreases as data store buffers increase in size.

Figure 4b illustrates the relationship between data cache miss rates and data store buffer sizes. In contrast to the relationship shown in Figure 4a, data cache miss rates are lowest when the corresponding data store buffer sizes are small. Data cache miss rates increase as data store buffers increase in size.

Figure 4c illustrates the relationship between overall cache-miss rates (data cache and instruction cache miss rates combined) and data store buffer sizes. A data store buffer size closest to the bottom 430 of the curve is optimal. The software decoder implementation of the invention selects a buffer size in accordance with a low overall cache miss-rate. Reliance upon this technique leads to performance improvements of several orders of magnitude for varying decoder architectures.

TABLE 2

An Example of a Video Decoding Module

```

5      // Cache Management Procedure

      read bitstream data           // encoded image data

      let N = 4                     // programmable number of macroblocks
10
      set each Data Store Buffer Size = 2*N    // data store buffer size is some
                                              // function of "N"
      let MAX_MB = 16               // maximum number of macroblocks in
                                              // encoded image data
15
      let MB = 0                    // macroblock initialization

      do while (MB < MAX_MB){

20          let i = MB
          for N number of macroblocks{
              VLD_IQ(i)              // Variable Length Decoder ("VLD")
                                      // decodes N
              i = i + 1              // macroblocks worth of bitstream data
25          }
          let i = MB
          for N number of macroblocks{
              IDCT(i)                // Inverse Discrete Cosine Transform
                                      // ("IDCT")
30          i = i + 1                // transforms N macroblocks of
                                      // coefficients into
                                      // DFD information
          }

```

let I = MB

for N number of macroblocks{

 MC(i) // Motion Compensation ("MC") Unit

 I = i + 1 // reconstructs N macroblocks of DFD

5 } // information

let MB = MB + N

}

10 // Analysis Procedure

Create a set of Performance Results based on

Cache Miss rates in relation to the selected Data Store Buffer Size

15 // Report Procedure

Report at least a subset of the Performance Results

20 The operation of the exemplary implementation of a Video Decoding Module shown in Table 2 is written such that each of the key computational software components of the decoding pipeline work over a varying number of macroblocks. The exemplary implementation is explained in the context of a standards-based software decoding pipeline as shown in Figure 2.

25 The exemplary implementation begins by reading the bitstream data. The bitstream data is encoded digital image data and can be accessed by the video decoding module in a variety of ways, such as reading the data from memory, or receiving the data in real-time from transmission media (e.g., satellite, over-the-air, and CATV). Next, the macroblock resolution, "N," is set to equal to some number of
30 macroblocks. The value of "N" can be set in a number of ways, e.g., at compile time, read from a data file, or even determined dynamically by applying an objective criterion to the selection process. Recall that prior art systems usually process a single

Alternately, step 560 may be implemented to measure a predetermined number of iterations.

5 If the condition tested at step 560 is not satisfied, processing returns to step 520. If the condition tested at step 560 is satisfied, the data store buffers are then set (step 570). As previously indicated, the buffer size adjuster module 152 may be used to set the data store buffer size. This data store buffer size may be fixed (static) during subsequent processing. Alternately, the data store buffer size may be dynamically changed by intermittently or continuously processing the encoded image data in accordance with the processing steps of Figure 5.

10 Those skilled in the art will appreciate that the processing steps of Figure 5 can be modified for a hardware decoder. In such a case, the step of determining cache performance results (step 540) would be substituted with the step of determining memory performance results, while the step of determining optimal cache performance behavior (step 560) would be substituted with the step of determining optimal memory
15 performance behavior.

Observe that the value "N" can also represent the amount that the code and data should be "unrolled" in the video decoding pipeline. The implementation of the video decoder with scalable buffering, as disclosed above, can alternatively be used by a video decoder architecture designer to determine optimal hardware buffer sizes.

20 The present invention can be implemented as a computer program product that includes a computer program mechanism embedded in a computer readable storage medium. For instance, the computer program product could contain the program modules shown in Fig. 5. These program modules may be stored on a CD-ROM, magnetic disk storage product, or any other computer readable data or program storage
25 product. The software modules in the computer program product may also be distributed electronically, via the Internet or otherwise, by transmission of a computer data signal (in which the software modules are embedded) on a carrier wave.

The foregoing description, for purposes of explanation, used specific nomenclature to provide a thorough understanding of the invention. However, it will
30 be apparent to one skilled in the art that the specific details are not required in order to practice the invention. In other instances, well known circuits and devices are shown in block diagram form in order to avoid unnecessary distraction from the underlying

